

Def.: A "stochastic learning algorithm" is an assignment  $S \mapsto \mu_S$ , where  $S \in (X \times Y)^{\mathbb{N}}$  and  $\mu_S$  is a prob. measure over a subset  $\mathcal{F} \subseteq \mathcal{Y}^X$  (with suitable  $\sigma$ -algebra, usually Borel).

We will write  $R(\mu_S) := \mathbb{E}_{h \sim \mu_S} [R(h)]$ ,  $\hat{R}_S(\mu_S) := \mathbb{E}_{h \sim \mu_S} [\hat{R}_S(h)]$

Exp.: Let  $r$  be a random variable with prob. d.s  $r \mapsto q(r)$  and  $A_r: S \mapsto h_S$  a learning alg. for every  $r$ . Then we can define

$$\mu_S(H) := \int q(r) \mathbb{1}_{A_r(S) \in H} dr \quad \text{for any measurable } H \in \mathcal{F}.$$

Then  $R(\mu_S) = \mathbb{E}_r [R(h_S)]$ ,  $h_S := A_r(S)$ .

Def.: A stochastic learning alg. is "uniformly stable" with rate  $\epsilon: \mathbb{N} \rightarrow \mathbb{R}$ , if  $\forall n \in \mathbb{N} \forall S, S' \in (X \times Y)^{\mathbb{N}} \forall (x, y) \in (X \times Y)$ :

$$\text{Ham}(S, S') = 1 \Rightarrow \left| \mathbb{E} L(y, h_S(x)) - \mathbb{E} L(y, h_{S'}(x)) \right| \leq \epsilon(n)$$

where the expectation refers to the stochastic component of the learning alg. (e.g.  $\mathbb{E} = \mathbb{E}_r$  in case  $h_S = A_r(S)$ ).

Implications of stability carry over to stochastic alg.s. For instance, one easily obtains on-average generalization in the form

$$\mathbb{E}_S [R(\mu_S) - \hat{R}(\mu_S)] \leq \epsilon(n)$$

The notion of differential privacy was introduced by C. Dwork (~2006) in the context of data-base analysis w.r.t. privacy. It turns out to imply stability:

Def.: A stochastic alg.  $S \mapsto \mu_S$  is " $\epsilon$ -differentially private" if for all measurable  $H \subseteq \mathcal{F}$  and all  $S, S'$  with  $\text{Ham}(S, S') = 1$ , we have:

$$\mu_S[H \in \mathcal{H}] \leq e^\epsilon \mu_{S'}[H \in \mathcal{H}]$$

Cor.: If a learning alg. is  $\epsilon$ -d.p., then it is  $2(e^\epsilon - 1)$  uniformly stable.

proof:

$$\left| \mathbb{E} L(y, h_S(x)) - \mathbb{E} L(y, h_{S'}(x)) \right| = \left| \int_{\mathcal{H}} L(y, h(x)) [\mu_S(h) - \mu_{S'}(h)] dh \right|$$

Hölder

$$\downarrow$$

$$\leq \underbrace{\sup_{h, x, y} |L(y, h(x))|}_{\leq 1} \int |\mu_S(h) - \mu_{S'}(h)| dh$$

$$\leq 2 \int \mu_S(h) - \mu_{S'}(h) dh$$

$$h: \mu_S(h) \geq \mu_{S'}(h)$$

$$\leq 2 \int \mu_{S'}(h) \left( \frac{\mu_S(h)}{\mu_{S'}(h)} - 1 \right) dh$$

$$\leq 2 (e^\epsilon - 1)$$

□